

Determining Health Utilities through Data Mining of Social Media

CL Thompson, Josh Introne, and Clint Young

(Dated: August 16, 2016)

Abstract

‘Health utilities’ measure patient preferences for perfect health compared to specific unhealthy states, such as asthma, a fractured hip, or colon cancer. When integrated over time, these estimations are called quality adjusted life years (QALYs). Until now, characterizing health utilities (HUs) required detailed patient interviews or written surveys. While reliable and specific, this data remained costly due to efforts to locate, enlist and coordinate participants. Thus the scope, context and temporality of diseases examined has remained limited.

Now that more than a billion people use social media, we propose a novel strategy: use natural language processing to analyze public online conversations for signals of the *severity* of medical conditions and correlate these to known HUs using machine learning. In this work, we filter a dataset that originally contained 2 billion tweets for relevant content on 60 diseases. Using this data, our algorithm successfully distinguished mild from severe diseases, which had previously been categorized only by traditional techniques. This represents progress towards two related applications: first, predicting HUs where such information is nonexistent; and second, (where rich HU data already exists) estimating temporal or geographic patterns of disease severity through data mining.

I. FULL TEXT:

The game theorist John Von Neumann and his collaborator, Morgenstern, designed one of the earliest measures of health utility (HU).¹ Their method, called the Standard Gamble, quantifies quality-of-life by first asking patients to make a hypothetical, important decision. Ultimately, the information contained in that decision is converted to a number ranging between 0 and 1. To start, a patient imagines that researchers have developed a potent drug that can cure the patient's disease. However the pill has a terrible side effect in some patients— instant death. So the patient is asked to decide on the maximum acceptable risk (m), expressed as a probability, that would allow proceeding with the therapy. A HU (u) is defined as $1 - m = u$.

Patients with a lower quality of life (lower u) have a stronger desire for health and accommodate more risk (higher m). Importantly, health utilities are not disease-specific, so diverse health conditions and outcomes can be compared. For example, we can say that epilepsy is worse than asthma because, on average, epilepsy patients accept a higher risk of death when faced with the dilemma in the Standard Gamble. In contrast, there is no clear hierarchy for *disease-specific* endpoints such as seizure frequency or pulmonary function.

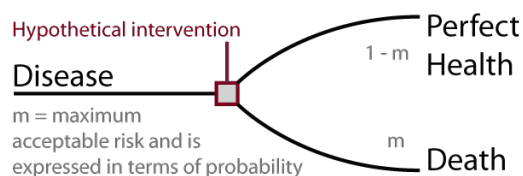


FIG. 1: Standard Gamble

The integral of a HU, with respect to time, is called a Quality Adjusted Life Year (QALY). This is a metric widely used in cost-utility analyses which assess relative economic value of various health interventions. Therefore, thousands of published studies (including clinical guidelines and CDC assessments of public health) rely on estimates of disease severity in the form of HUs.²

For decades, questionnaires or interviews like the Standard Gamble have served to elicit health utilities. Other examples in this class include: Time Trade Off and the Health Utilities Index. An alternative involves using algorithms that map clinical data to HUs. Generally, for a given health condition, clinical data and HU data must be collected simultaneously; then a regression model is built such that the objective data predicts the HU. For example, epilepsy patients with more frequent seizures may report lower quality-of-life, (although some evidence confounds this assumption).³ These mappings are generally limited to their disease context. For example, knowledge about the relationship of seizure frequency to HU does not transfer to a mapping of pulmonary function to HU.

Most recently, investigators have utilized natural language processing to filter Twitter and other social networks for novel indices of disease severity. Generally, HUs have not been utilized as training labels for these metrics, and algorithms apply to specific conditions only, such as toothaches⁴ and depression.⁵ Despite vast repositories of patient generated data, there is no general data mining approach that can estimate a HU for any given disease.

Substantial progress has been made by Parimbelli *et al.*⁶ who have used sentiment analysis (SA) of online health messaging boards to derive HUs. SA uses natural language processing (NLP) to extract emotional information from online messages; corporations have used it to monitor brand appeal and word-of-mouth recommendations by consumers. Parimbelli *et al.* use NLP to identify relevant content within five wellness domains (e.g. mobility, presence of pain, *etc.*). Then, using SA, a negativity score is created for each domain, providing weights for its HU estimate. The application computes a final score by weighting these HUs with traditionally derived HUs. So far the system has been validated on a scenario related to atrial fibrillation.

Like Parimbelli, we have an ambitious scope: we train an algorithm to estimate a HU for potentially any disease— not just a mapping that transforms disease-specific measurements into

HUs. In contrast to Parimbelli, we use alternate natural language features to estimate HUs. We also highlight the performance of our method, comparing its estimates to held out data on 20 traditionally-derived HUs.

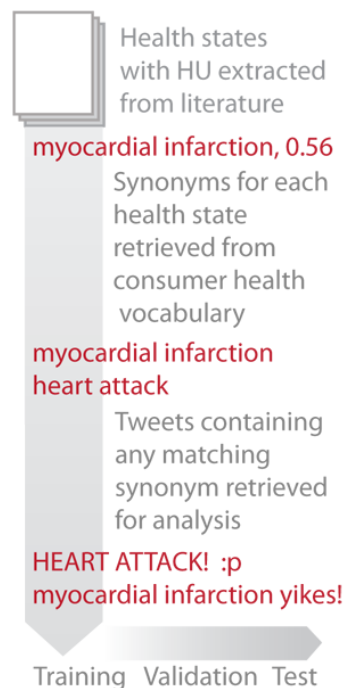


FIG. 2: Process

We performed the search for related messages by using the exact term cited in the study, and then expanded our query to include all matching synonyms as provided by a consumer health vocabulary containing validated mappings between expert and layperson terminology.⁸ We defined the entire corpora of social media messages related to a single disease, along with its class label, as ‘an example’.

Our social media content came from Twitter, a ‘microblogging’ site that limits user messages (called, ‘tweets’) to 140 characters. Challenges of using Twitter include its unmoderated content, unbounded by topic or context. Yet, users generate more than 500 million tweets every day and Twitter data has been utilized for hundreds of epidemiological studies. Data has been published showing that people share highly personal information.⁹ We obtained a dataset of 11 million health-related tweets, previously filtered from a dataset of two billion tweets.¹⁰

Finally we followed a standard methodology for machine learning experiments: each illness was randomly assigned to training, validation, and test sets. All sets had roughly the same proportion of mild and severe illnesses.

To develop a candidate set of features, we took two complementary approaches. For the first approach, we inspected a small sample of tweets. Scrolling through the data revealed that topics related to mild health problems often included a frivolous style such as usage of emoticons, while severe diseases engendered a more formal, grave tone that included conventional punctuation and syntax. Therefore, our candidate features became emoticons and patterns related to punctuation or capitalization.

For the second approach, we relied on large scale screening to select single words (unigrams, hereafter). As a convenient source, we took a sentiment dictionary containing 2477 unigrams used for sentiment analysis (SA),¹¹ then found the frequency of these unigrams in our social media corpus. Selecting the most frequent 50 unigrams, we further refined this set to include 16 unigrams

To test the idea that HUs can be derived from online health conversations, we attempted to build an algorithm to classify mild versus severe health conditions– the most essential dichotomy. We began by compiling a list of diseases with HUs previously determined by well-designed, highly-cited studies utilizing the traditional approach of interviews or surveys.

We preferred highly cited studies with the most scope in order to reduce the variation that inevitably arises between HU studies. Some studies undertake broad data collection efforts, apply a uniform methodology and report HUs for more than 100 diseases. Also, we required sufficient specificity with respect to the health conditions studied. Some HU studies aggregate conditions by International Classification of Disease code. Unfortunately, this lossy process results in terms that are too broad or vague for our purposes, such as “esophageal problem.”⁷

For selected studies, we extracted health conditions with their HUs as ‘gold-standard’ labels. The conditions with the lowest HU became the ‘severe’ set, while the upper half became the ‘mild’ conditions, although we culled a margin of middle range of HUs (see table 2 below). Then for each condition in each class, we found related social

	Regular Expression	Notes on Matching Pattern	Example: a tweet with a matching pattern highlighted in red.
1	[A-Z][!]	Capital followed by exclamation.	tachycardia OMG!!
2	[A-Z][A-Z][A-Z][A-Z]	4 consecutive capitalized letters.	...really being hit in the throat by hay fever. ANNOYING
3	[:=][(P/Oo]	Emoticon such as “ : / ” or “ : O ”	Lay off on the coffee. caffeine can trigger gout according to the doctors. :P
4	[]([Rr][Ii][Ss][Kk]	Word forms of “risk” or “RISK”	Are you at risk for prostate cancer? Early detection is key.
5	[]([Hh][Aa][Tt][Ee][]	“ hate ” capitalization insensitive	I just woke up in the emergency department. I hate having epilepsy...
6	[]([Ff][Ii][Gg][Hh][Tt][]	“fight ” capitalization insensitive	Please increase support for the Global Fund to fight AIDS, and tuberculosis.
7	[]([Ww][Oo][Rr][Ss][Ee][]	“ worse ” capitalization insensitive	Just caught myself stuttering kinda like but worse
8	[]([Dd][i][e]	Word forms of “ die ” or “ Die ”	The best thing about tachycardia is: I could die any moment. honest.
9	[]([Dd][e][a][t][h][]	“ death ” or “ Death ”	TB is a leading cause of death for HIV patients.

TABLE I: Features

that best separated the classes.

While we screened more than 60 variables for relevance, we selected only 24 features for training, which included 16 unigrams, three emoticon features, and five formality patterns. Using backward selection, we eliminated 15 features that had a minimal contribution to the classification accuracy on the validation set, leaving nine features in the final algorithm.

For our learning algorithms, we compared the following classifiers implemented in the Python library, Scikit-learn: logistic regression, a decision tree, a linear support vector machine, and a Random Forest. These algorithms represent widely known, general-purpose algorithms in machine learning.

The decision tree classifier (with default settings) performed the best of the four algorithms: its accuracy is 72% on the test set. The final algorithm includes nine features, including emoticons, unigrams, and patterns indicating formality of sentence construction. The algorithm tends to misclassify mild diseases as severe. Amid these misclassified examples, there are no obvious trends in terms of data characteristics, such as number of tweets in the test set, nor are there obvious clinical features linking these diseases, such as age of onset, duration, or gender distribution.

Thus we have shown modest accuracy despite several challenges: our social media content is wide-ranging, and users are only a subset of the population. We did not extensively filter tweets to narrow context, such as for self-experience versus unrelated content. Also, the volume of tweets

Training	HU	Validation	HU	Test	HU	Error
Lung cancer	0.39	Hydrocephalus	0.19	Blind	0.47	
Aneurysm	0.49	Heart disease	0.56	Emphysema	0.48	
Atherosclerosis	0.52	Prostate cancer	0.56	Dementia	0.54	
Angina	0.56	Tuberculosis	0.57	Hip fracture	0.60	
Breast cancer	0.64	Cerebral palsy	0.60	Brain cancer	0.60	
Cataracts	0.65	Liver cancer	0.61	Diabetes	0.62	
Constipation	0.66	Epilepsy	0.62	Kidney infection	0.63	x
Herniated disk	0.66	Renal failure	0.65	Deaf	0.65	
Arrhythmia	0.66	Gout	0.65	Arthritis	0.69	
Hernia	0.67	Tachycardia	0.65	Depression	0.70	
Bursitis	0.74	Anxiety ¹	0.68	Hearing impairment	0.74	
Bunion	0.76	Tendinitis	0.73	Asthma	0.77	
Callus	0.77	Thyroid disorder	0.74	Psoriasis	0.80	x
Scoliosis	0.77	Indigestion	0.74	Flat feet	0.80	
Murmur	0.79	Stutter	0.74	Hemorrhoids	0.84	x
Dry skin	0.80	Hay fever	0.87	Acne	0.89	x
Sinusitis	0.82	Allergic rhinitis	0.87	Allergic reaction	0.91	
Dermatitis	0.83	Color blind	0.89	Alopecia	0.91	
SEVERE VS. MILD		ADHD	0.92	Gastroenteritis	0.92	x
x = MISCLASSIFIED		Otitis media	0.96	Hiccups	0.99	

TABLE II: Health Conditions

on each health condition is quite variable, ranging from about 100 tweets to 200,000. Since we had to manually extract suitable HUs from the literature, we have a limited supply of examples for training and validation.

Algorithm	Validation	Test
Decision Tree	74	72
Random Forest	75	71
Support Vector Machine	75	65
Logistic Regression	75	65

TABLE III: Accuracy

Limitations of data mining include the following: as noted in similar studies of online patient generated data, periodic crowd-level events cause unpredictability.¹² Using a diverse set of predictive features, and filtering for specific contexts could mitigate this to some degree. Also, attention will have to be paid to who is posting messages to social media and how they advocate or otherwise act as a proxy for patients with a particular disease. Sometimes HUs are obtained within highly controlled situations, such as clinical drug trials; clearly data

mining cannot replace face-to-face elicitation in these circumstances.

Data mining addresses the many gaps in knowledge that exist (especially for the developing world) because of the expense of the traditional methodology. Our work represents a movement towards an information-age scaling of a clinical, research, and quality improvement task: gauging patient complaints related to health conditions and outcomes. In the business world, data mining of consumer sentiment is already widespread. Epidemiologists have already developed systems for tracking case counts using search engine data,¹² Twitter,¹³ or Wikipedia pageviews.¹⁴ Much of this work has focused on case counts but has lacked an assessment of case severity, and therefore misses an aspect of disease burden.

In summary, we have used 11 million health-related tweets and a learning algorithm to classify health conditions as severe or mild, a necessary step in the final goal of placing conditions on a continuum of severity. The power of data mining is that it can be applied to any disease with measurements occurring at any time interval. Unlike a survey, which has been the mainstay for more than 70 years, the marginal cost of analyzing an additional health condition is virtually zero. In particular, this has implications for the developing world, where text messaging and social media have grown increasingly prevalent and yet the assessment of disease burden remains challenging.

II. METHODS

Literature Selection: We used Pubmed and Google Scholar to identify established studies that defined health utilities for multiple, specific health conditions. Most of our HUs were obtained from a single study,⁷ sponsored by the National Center for Health Statistics (NCHS). It met our criteria since it included 130 conditions and (as an indicator of acceptance) had 268 citations as of April 23, 2016. This agency, charged with monitoring the health of the US population for the Department of Health and Human Services, conducted the research to validate a novel HU measurement instrument (called HALex). The authors of the NCHS study conclude that HALex has good face validity and high correlations to other canonical studies. A secondary goal was to generate a catalog of HUs useful to studies that must rely on HU information from secondary data sources due to the prohibitive expense of acquiring primary data. We bolstered our HU dataset with other studies cataloging HUs,¹⁵ including two containing pediatric health conditions^{16,17} in order to diversify the conditions studied. A few conditions were included that were not part of studies assessing multiple health conditions simultaneously.^{18–20} To mitigate errors related to methodological variations between studies, we preferred that these single-study conditions had clear class membership (very mild or very severe).

Selection of Concepts: Although the NCHS study reports HUs on many specific illnesses, it still contains HUs associated with unhelpful, broader diagnostic bins: examples include ‘other extremity paralysis’, ‘orthopedic impairment-other’, ‘absent bone/joint’, and ‘hand or finger impairment’. We used common experience and the medical expertise of one of the authors, CLT, a medical doctor, to filter non-specific terms. In about 30 cases, we could not use a selected concept since we had an insufficient number of tweets (arbitrarily, fewer than 100).

Selection of Tweets: We filtered our dataset of 11 million health-related tweets for content related to our research. First, we selected all the tweets containing our concept of interest, matching tweets with a capitalization-insensitive algorithm. A concept is comprised of multiple descriptions (synonyms). For each HU study, we found a summary table listing every health condition measured and its corresponding HU. We extracted the exact health description mentioned in the summary table, and then retrieved all of its descriptions and pertinent word forms as defined by a consumer health vocabulary mentioned above. If a tweet matched more than one concept, it was selected for each concept.

Privacy: The health-related tweets obtained for this study contained only the content of the tweet, but we eliminated all other metadata. Because the tweets were intentionally public, this study was exempt from the IRB process. The tweets in table 1 were also public, and have been further modified by us to obscure attribution to a specific user.

Algorithm Development: For our learning algorithms, we compared the classifiers as discussed in the body of this article and implemented in the Python library, Scikit-learn. We randomized the examples to training, validation and test sets using stratified randomization, so that each set contained a similar range of health utilities. As discussed, we sequentially explored the decision tree hyperparameters by simply assigning increasingly extreme values while monitoring the effect on the validation error. As discussed above, we selected candidate predictive features via two approaches: manually developing features that capture differences in formality and in the second approach, differentiating frequencies of unigrams from the affective lexicon. Nuances related to the capitalization of these unigrams change their predictive value.

III. REFERENCES

-
- ¹ Neumann, J. v. & Von Morgenstern, O. Theory of Games and Economic Behavior. *Princeton University Press* (1944).
 - ² Neumann, P. J. & Weinstein, M. C. Legislating against use of cost-effectiveness information. *New England Journal of Medicine* **363**, 1495–1497 (2010).
 - ³ Choi, H. *et al.* Seizure frequency and patient-centered outcome assessment in epilepsy. *Epilepsia* **55**, 1205–1212 (2014).
 - ⁴ Heavilin, N., Gerbert, B., Page, J. E. & Gibbs, J. L. Public health surveillance of dental pain via Twitter. *Journal of dental research* **90**, 1047–1051 (2011).

- ⁵ De Choudhury, M., Counts, S. & Horvitz, E. Social media as a measurement tool of depression in populations. In *Proceedings of the 5th Annual ACM Web Science Conference*, 47–56 (ACM, 2013).
- ⁶ Parimbelli, E. *et al.* Use of Patient Generated Data from Social Media and Collaborative Filtering for Preferences Elicitation in Shared Decision Making. In *2014 AAAI Fall Symposium Series* (2014).
- ⁷ Gold, M. R., Franks, P., McCoy, K. I. & Fryback, D. G. Toward consistency in cost-utility analyses: using national measures to create condition-specific values. *Medical care* **36**, 778–792 (1998).
- ⁸ Zeng-Treitler, Q., Goryachev, S., Kim, H., Keselman, A. & Rosendale, D. Making texts in electronic health records comprehensible to consumers: a prototype translator. In *AMIA*, 846–50 (2007).
- ⁹ Lee, J. L., DeCamp, M., Dredze, M., Chisolm, M. S. & Berger, Z. D. What are health-related users tweeting? A qualitative content analysis of health-related users and their messages on twitter. *Journal of medical Internet research* **16**, e237 (2014).
- ¹⁰ Paul, M. J. & Dredze, M. Discovering health topics in social media using topic models. *PLoS One* **9**, e103408 (2014).
- ¹¹ Nielsen, F. A. A new ANEW: Evaluation of a word list for sentiment analysis in microblogs. *arXiv preprint arXiv:1103.2903* (2011).
- ¹² Ginsberg, J. *et al.* Detecting influenza epidemics using search engine query data. *Nature* **457**, 1012–1014 (2009).
- ¹³ Lamb, A., Paul, M. J. & Dredze, M. Separating Fact from Fear: Tracking Flu Infections on Twitter. In *HLT-NAACL*, 789–795 (2013).
- ¹⁴ Generous, N., Fairchild, G., Deshpande, A., Del Valle, S. Y. & Friedhorsky, R. Global disease monitoring and forecasting with Wikipedia. *PLoS Comput Biol* **10**, e1003892 (2014).
- ¹⁵ Fryback, D. G. *et al.* The beaver dam health outcomes study initial catalog of health-state quality factors. *Medical Decision Making* **13**, 89–102 (1993).
- ¹⁶ Carroll, A. E. & Downs, S. M. Improving decision analyses: parent preferences (utility values) for pediatric health outcomes. *The Journal of pediatrics* **155**, 21–25 (2009).
- ¹⁷ Petrou, S. & Kupek, E. Estimating preference-based health utilities index mark 3 utility scores for childhood conditions in England and Scotland. *Medical decision making* **29**, 291–303 (2009).
- ¹⁸ Bilgi, . *et al.* Psychiatric symptomatology and health-related quality of life in children and adolescents with alopecia areata. *Journal of the European Academy of Dermatology and Venereology* **28**, 1463–1468 (2014).
- ¹⁹ Fluchel, M. *et al.* Self and proxy-reported health status and health-related quality of life in survivors of childhood cancer in Uruguay. *Pediatric blood & cancer* **50**, 838–843 (2008).
- ²⁰ Roux, C. *et al.* Burden of non-hip, non-vertebral fractures on quality of life in postmenopausal women. *Osteoporosis International* **23**, 2863–2871 (2012).

IV. END NOTES

Supplementary Information Table comprised of all health conditions used in the study and for each: number of tweets, HU value from literature, source of HU value, partition (training, validation, or test).

Acknowledgements The authors are grateful to Michael J. Paul and Mark Dredze for sharing Tweets, Reid Friedhorsky for sharing feedback, and Rebecca Anthony for her editing, critiques, and encouragement.

Author Contributions: C.L.T. conceptualized the study. J.I. provided feedback on the study design. C.L.T. and C.Y. performed the analyses. C.L.T. prepared the manuscript, with extensive

editing by J.I. and C.Y.

Author Information: The authors are affiliated with Michigan State University, (East Lansing, MI 48824, USA), from the following departments: Chris Thompson, MD, MHI is from pediatrics, Clint Young, PhD is from the department of Physics, Astronomy and the National Superconducting Cyclotron Laboratory. Josh Introne, PhD is from Media and Information. The authors declare no competing financial interests. The corresponding author is Christopher L. Thompson.